

A decorative graphic on the left side of the slide consists of a grid of colored squares. The grid is 4 squares wide and 4 squares high, with the bottom-right square missing. The colors of the squares are: Row 1: teal, orange, brown, brown; Row 2: orange, brown, light brown, light brown; Row 3: orange, teal, light brown, light brown; Row 4: light brown, orange, orange, brown.

Градиентные методы обучения

Повторение

Функция потерь для регрессии

- Еще один вариант – средняя абсолютная ошибка (mean absolute error, MAE)

$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l |a(x_i) - y_i|$$

- Слабее штрафует за выбросы
- Зануляет признаки перед незначущими признаками

Линейная регрессия в векторном виде

$$a(x) = \langle w, x \rangle$$

- Среднеквадратичная ошибка и задача обучения:

$$\frac{1}{l} \sum_{i=1}^l (\langle w, x_i \rangle - y_i)^2 \rightarrow \min_w$$

Модель линейной регрессии

Среднеквадратичная ошибка и задача обучения:

$$\frac{1}{l} \sum_{i=1}^l (\langle w, x_i \rangle - y_i)^2 \rightarrow \min_w$$

Отклонения прогнозов от ответов:

$$Xw - y = \begin{pmatrix} \langle w, x_1 \rangle - y_1 \\ \vdots \\ \langle w, x_\ell \rangle - y_\ell \end{pmatrix}$$

Вычисление ошибки

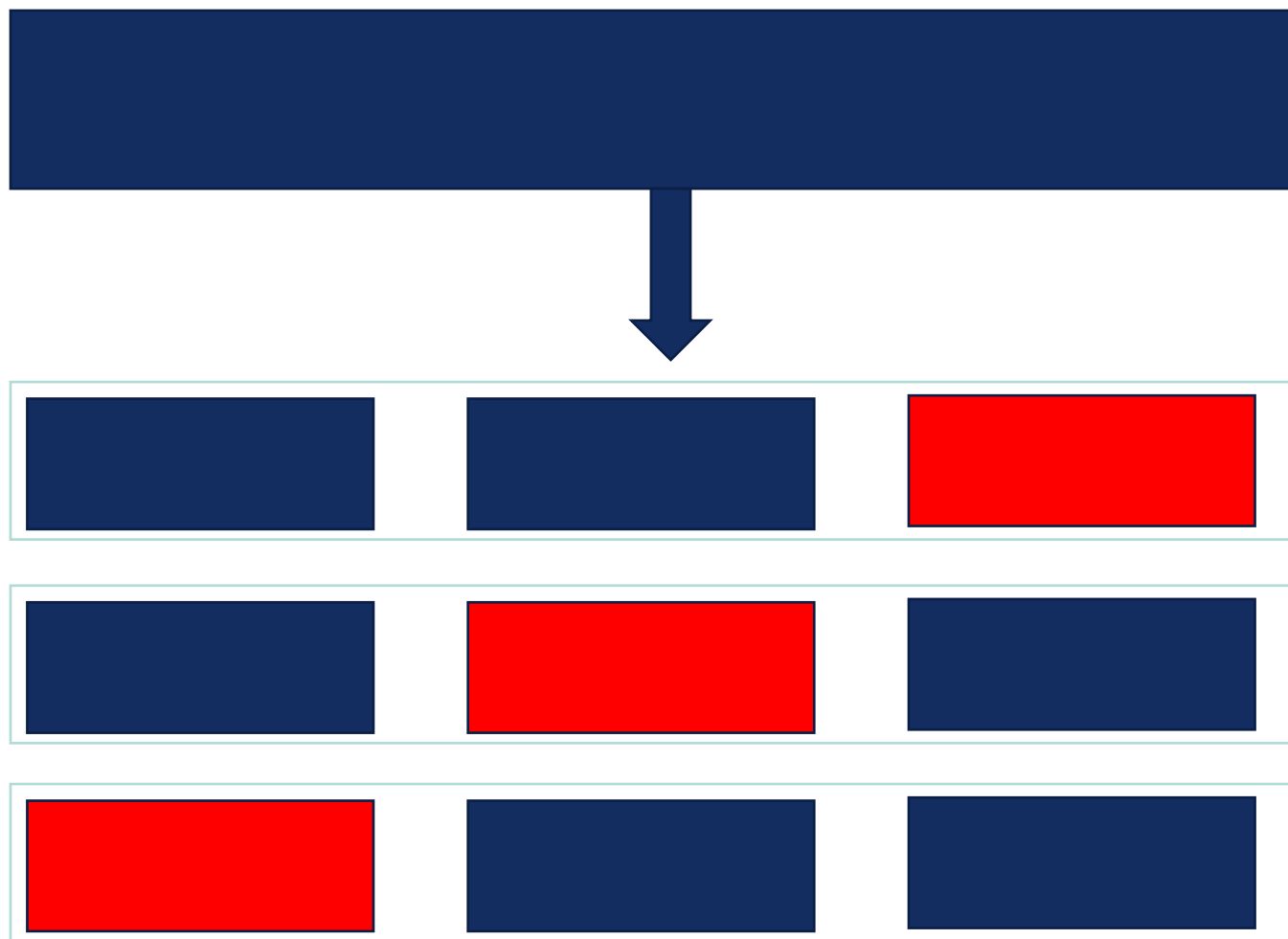
Отклонения прогнозов от ответов:

$$||z|| = \sqrt{\sum_{j=1}^n z_j^2}$$

Среднеквадратичная ошибка:

$$\frac{1}{l} ||Xw - y||^2 = \frac{1}{l} \sum_{i=1}^l (< w, x_i > - y_i)^2$$

Кросс-валидация



Признаки переобученной модели

- Большая разница в качестве на тренировочных и тестовых данных(модель подгоняется под тренировочные данные и не может найти истинная зависимость)
- Большие значения параметров (весов) w_j модели

Регуляризаторы

- $\|z\|_2 = \sqrt{\sum_{j=1}^d z_j^2}$ — L_2 — норма
- $\|z\|_1 = \sum_{j=1}^d |z_j|$ — L_1 — норма

Гиперпараметры

- Гиперпараметр алгоритма — то, что задается вручную
- Пример: коэффициент регуляризации
- Параметр алгоритма — то, что определяется моделью
- Пример: веса регрессии
- Как подбирать гиперпараметры алгоритма?
- Нельзя подбирать по обучающей выборке — это приведет к переобучению
- Нужно использовать дополнительные данные (валидация)

Чуть больше терминов

- После подбора всех гиперпараметров стоит проверить на совсем новых данных, что модель работает
- Обучающая выборка — построение модели
- Валидационная выборка — подбор гиперпараметров модели
- Тестовая выборка — финальная оценка качества модели

Масштабирование признаков

- Отмасштабируем j-й признак
- Вычисляем среднее и стандартное отклонение признака на обучающей выборке:

$$\mu_j = \frac{1}{l} \sum_{i=1}^l x_i^j$$
$$\sigma_j = \sqrt{\frac{1}{l} \sum_{i=1}^l (x_i^j - \mu_j)^2}$$

Среднеквадратичная ошибка

- MSE для линейной регрессии:

$$\frac{1}{l} \|Xw - y\|^2 \rightarrow \min_w$$

$$Q(w_1, \dots, w_d) = \frac{1}{l} \sum_{i=1}^l (w_1 x_{1i} + \dots + w_d x_{di} - y_i)^2$$

Можно посчитать градиент MSE:

- $\nabla \frac{1}{l} \|Xw - y\|^2 = \frac{2}{l} X^T (Xw - y)$

Приравниваем нулю и решаем систему линейных уравнений:

- $w = (X^T X)^{-1} X^T y$

Аналитическое решение

$$w = (X^T X)^{-1} X^T y$$

- Если матрица $(X^T X)$ вырожденная, то будут проблемы
- Даже если она почти вырожденная, все равно будут проблемы
- Если признаков много, то придется долго ждать
- Обращение матрицы – сложная операция ($O(N^3)$ от числа признаков
- Если заменить среднеквадратичный функционал на другой, то скорее всего не найдем аналитическое решение

Регуляризация

- Регуляризованный функционал

$$\frac{1}{l} \sum_{i=1}^l (\langle w, x_i \rangle - y_i)^2 + \lambda \|w\|^2 \rightarrow \min_w$$

- Аналитическое решение:

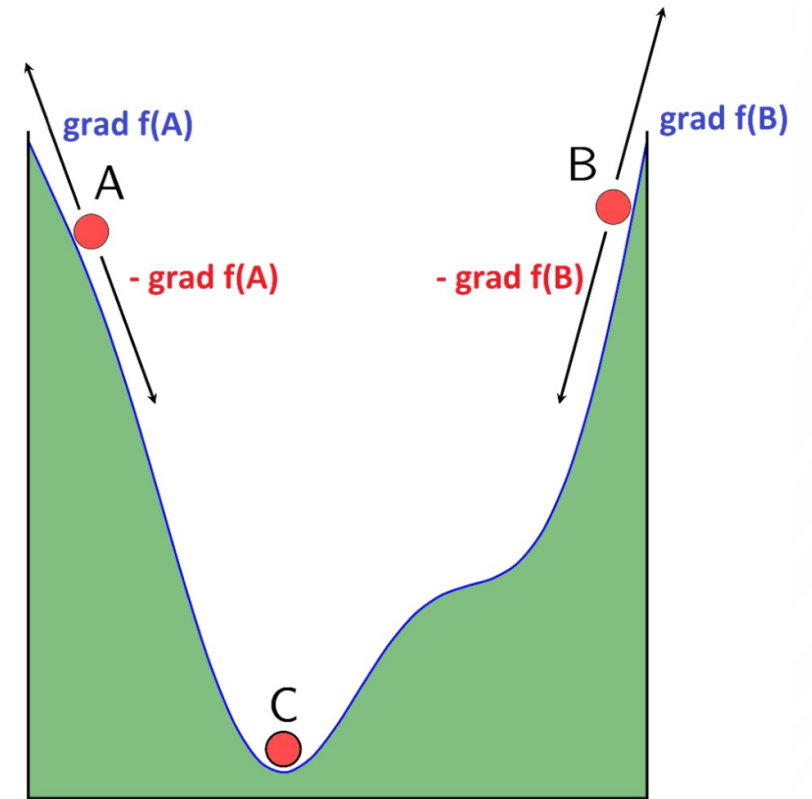
$$w = (X^T X + \lambda I)^{-1} X^T y$$

- Гребневая регрессия (Ridge regression)

Теорема о градиенте

Теорема. Градиент – это вектор, в направлении которого функция быстрее всего растет.

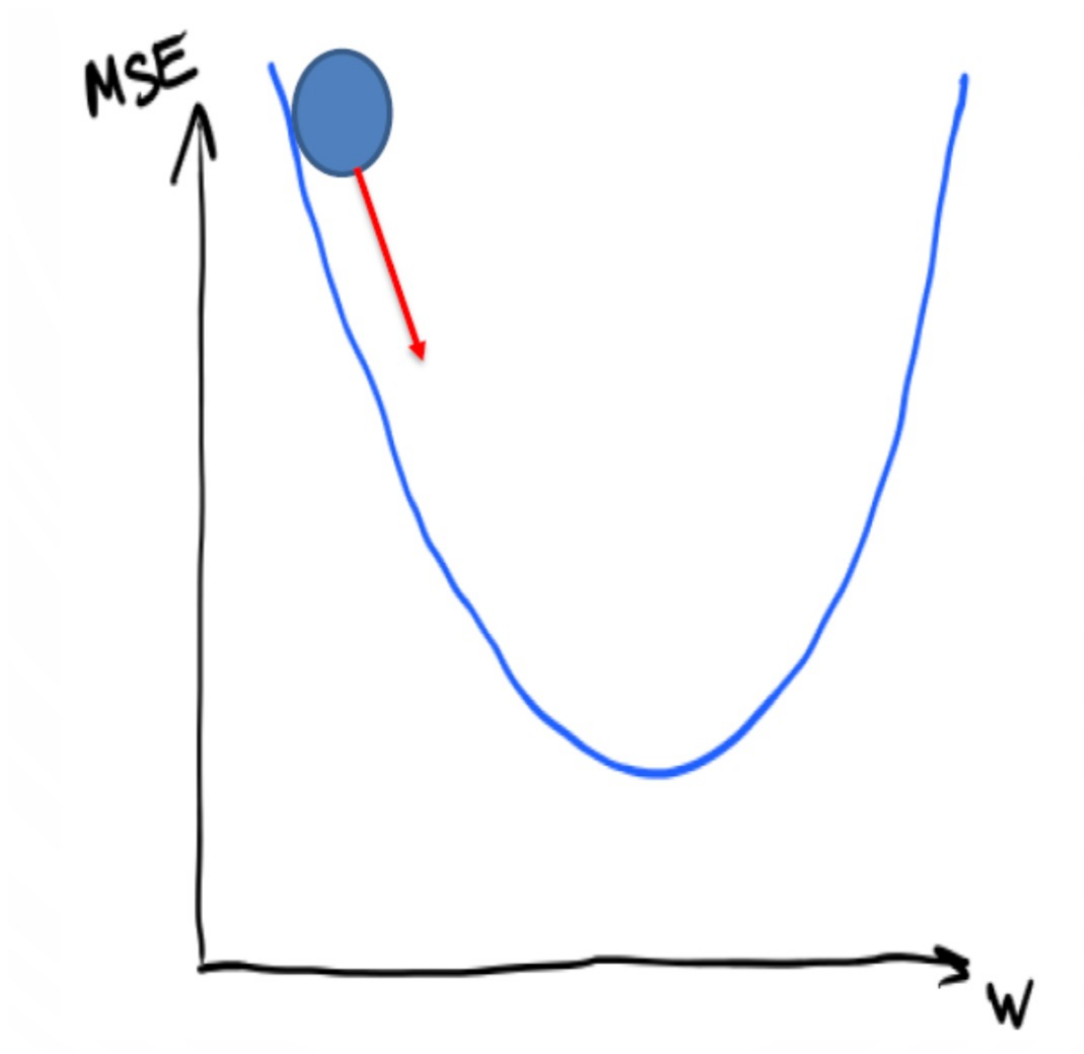
Антиградиент (вектор, противоположный градиенту) – Вектор, в направлении которого функция быстрее всего убывает.



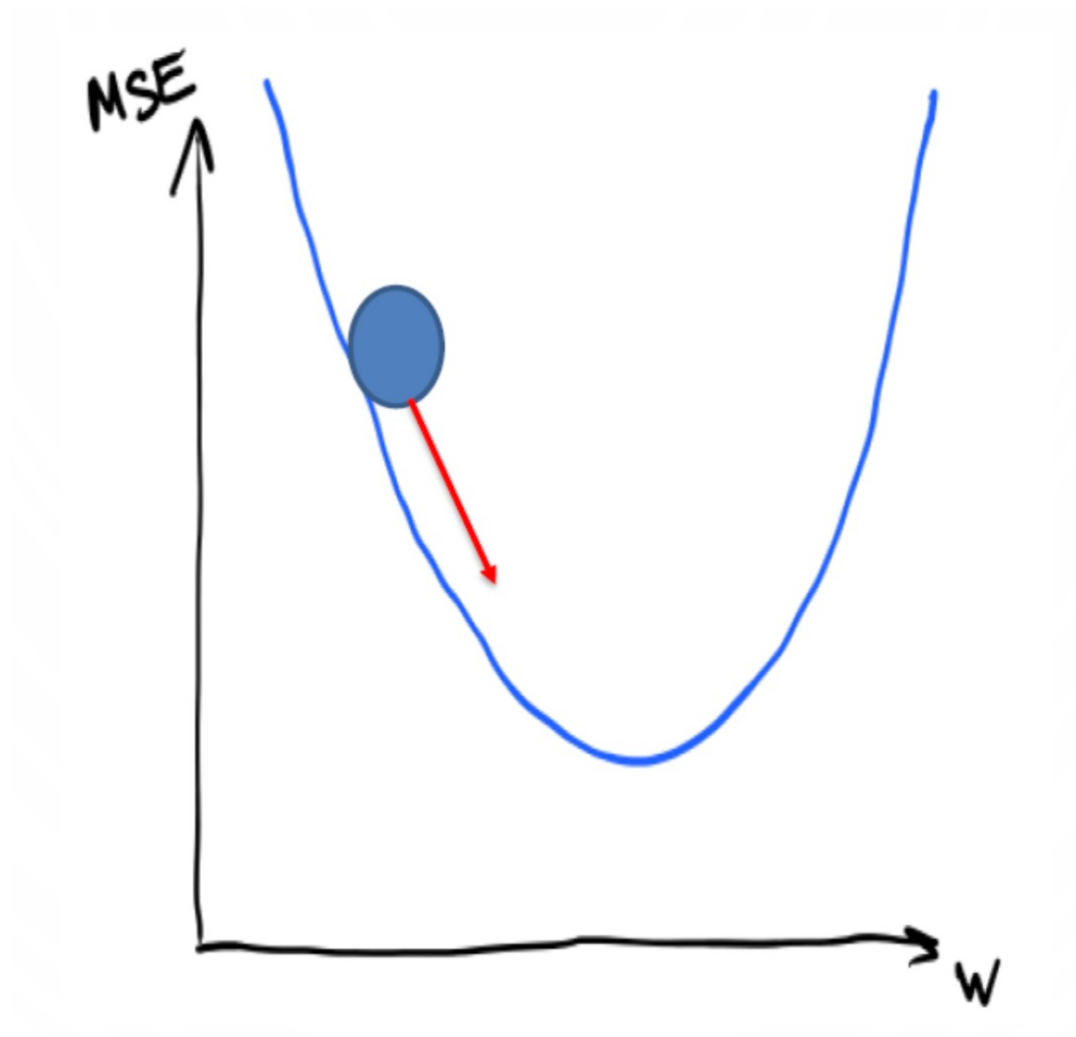
Метод градиентного спуска

- Наша задача при обучении модели – найти такие веса w , на которых достигается минимум функции ошибки.
- Грубо говоря, график MSE – парабола
- Идея метода градиентного спуска.
- На каждом шаге движемся в сторону антиградиента функции потерь!
- Вектор градиента функции потерь обозначают $grad\ Q$ или ∇Q

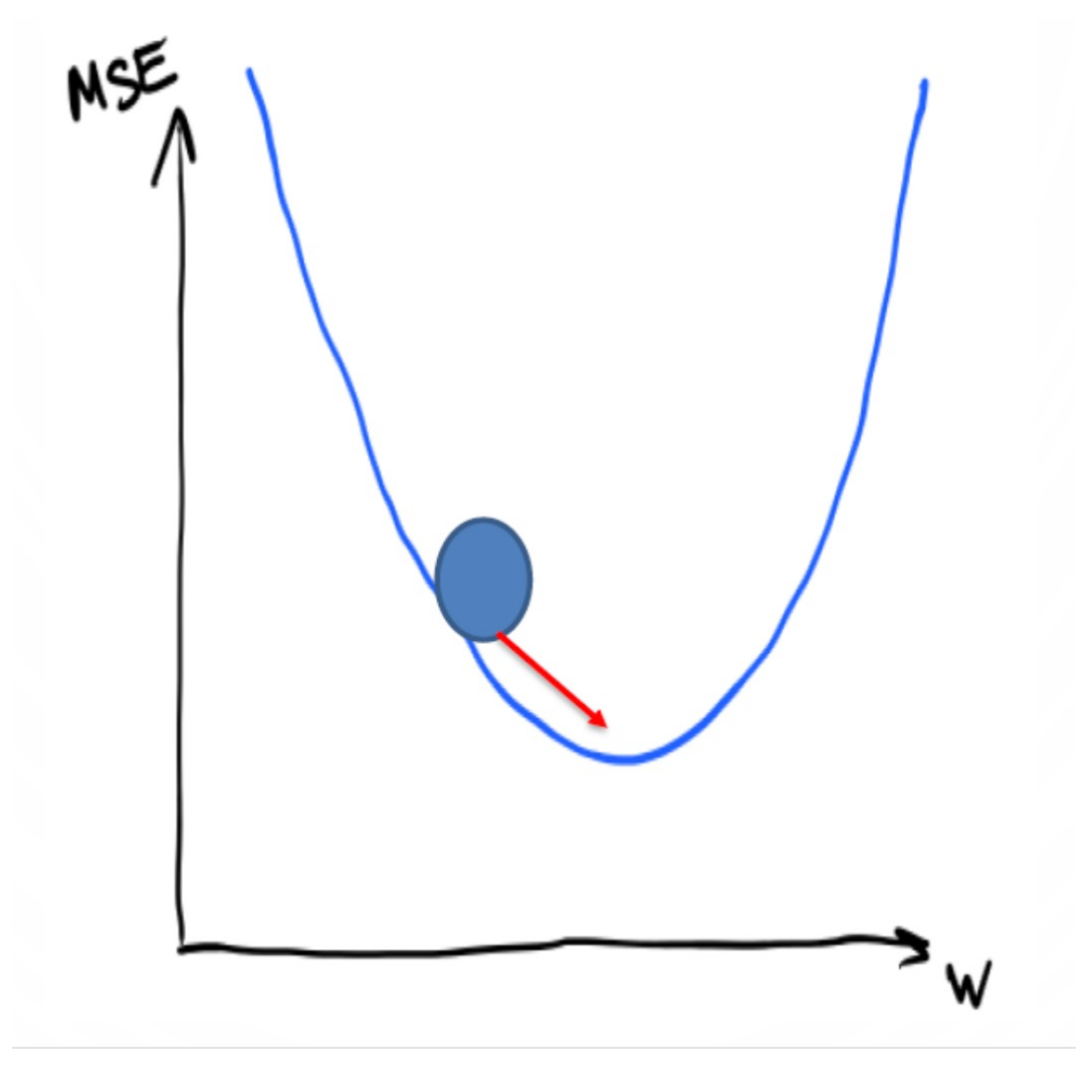
Метод градиентного спуска



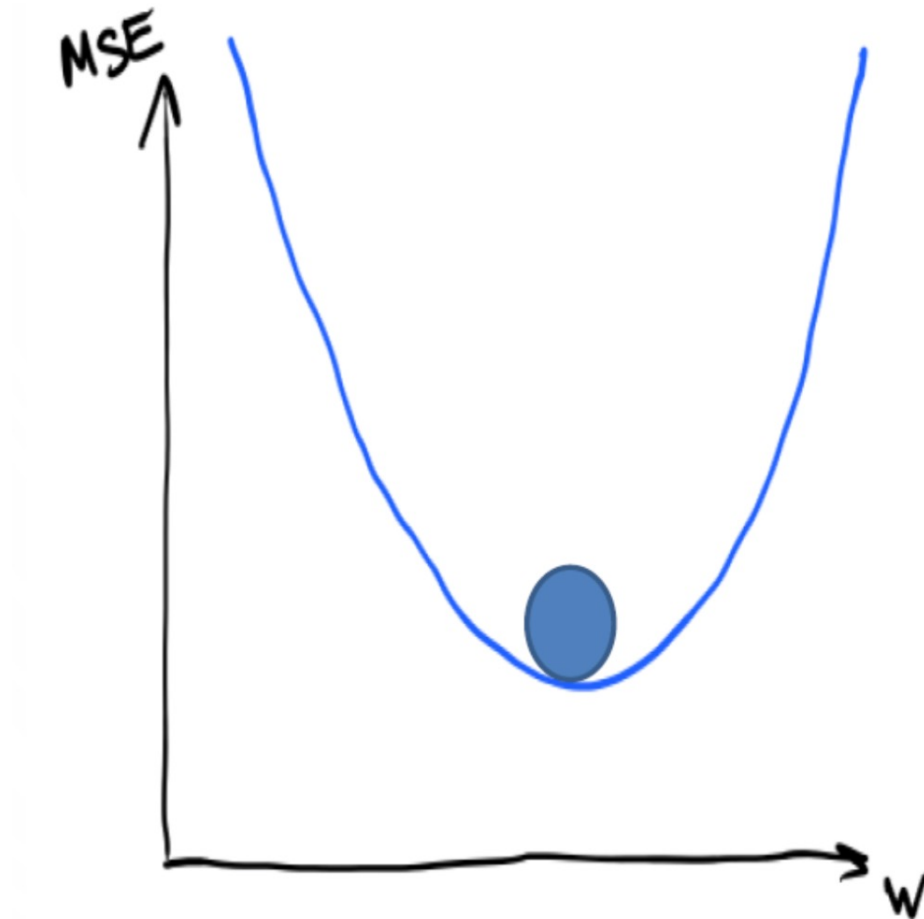
Метод градиентного спуска



Метод градиентного спуска

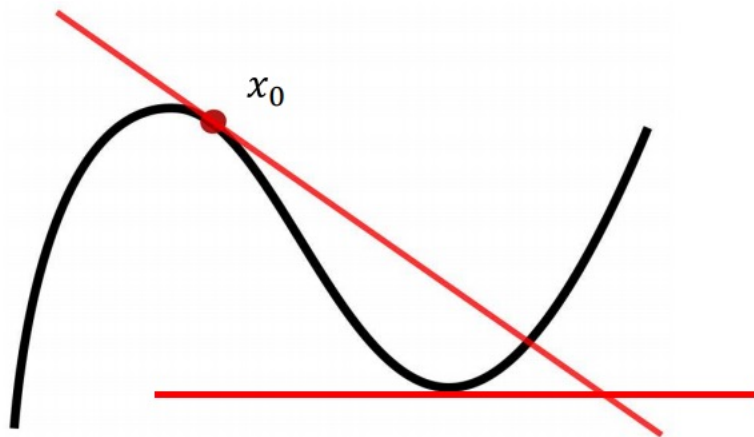


Метод градиентного спуска



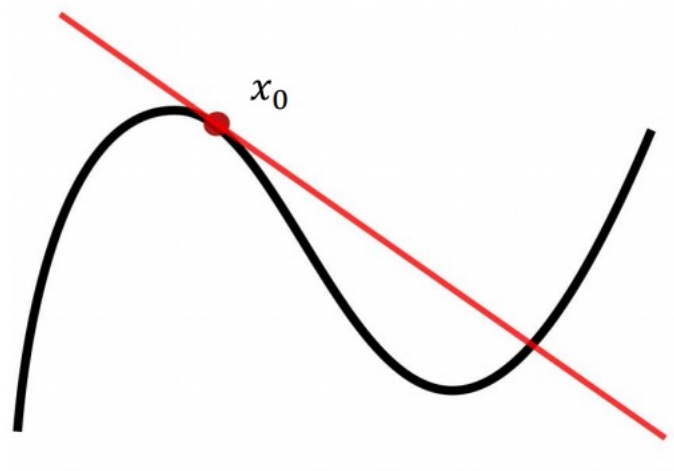
Производная

Если точка x_0 — экстремум и в ней существует производная, то $f'(x_0) = 0$



Производная

$$\lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} = f'(x_0)$$

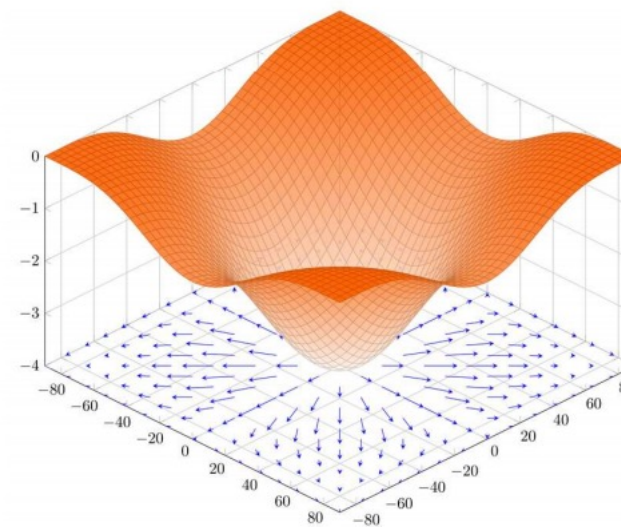


Градиент

- Градиент – вектор частных производных

$$\nabla f(x) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right)$$

- У градиента есть важное свойство!
- Зафиксируем точку x_0
- В какую сторону функция быстрее всего растёт?



Важное свойство

- Зафиксируем точку x_0
- В какую сторону функция быстрее всего растет?
- В направлении градиента!
- А быстрее всего убывает в сторону антиградиента

Условие экстремума

- Если точка x_0 - экстремум и в ней существует производная, то $\nabla f(x_0) = 0$
- Если функция выпуклая, то экстремум один
- MSE для линейной регрессии — выпуклая!

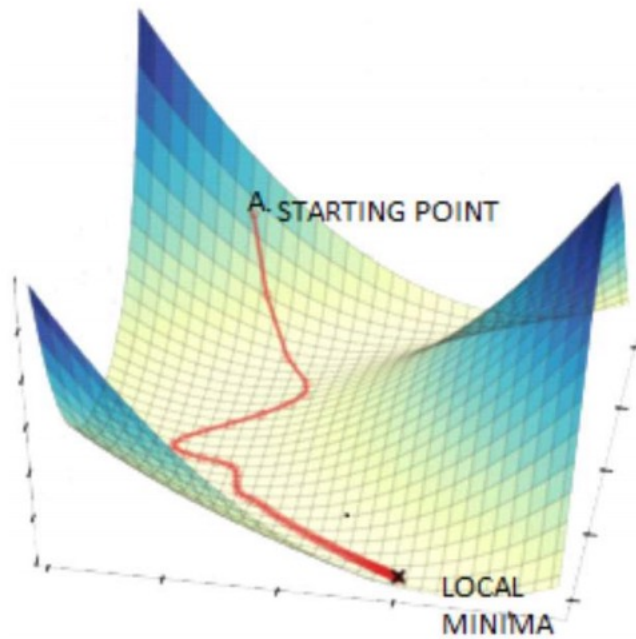
Градиентный спуск

- Если точка x_0 - экстремум и в ней существует производная, то $\nabla f(x_0) = 0$
- Если функция выпуклая, то экстремум один
- MSE для линейной регрессии — выпуклая!

Как это пригодится?



Как это пригодится?



Градиентный спуск

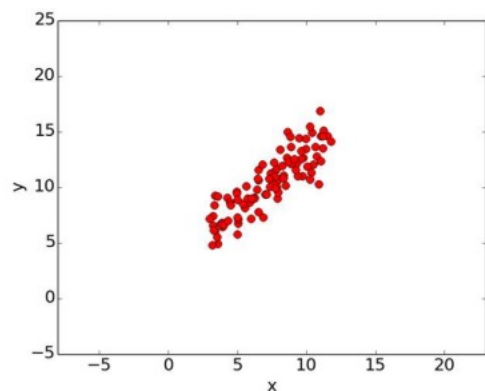
- Стартуем из случайной точки
- Сдвигаемся по антиградиенту
- Повторяем, пока не окажемся в точке минимума

Парная регрессия

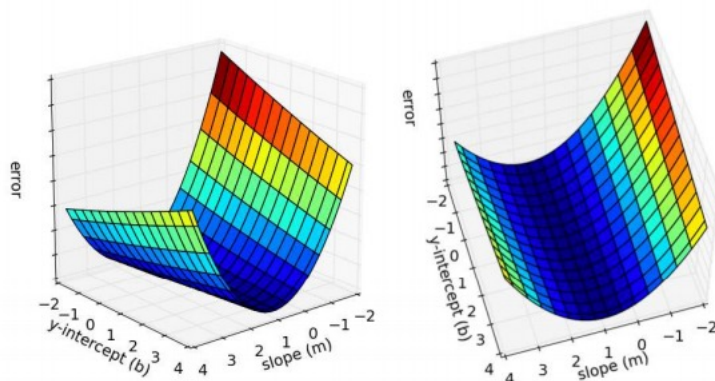
- Простейший случай: один признак
- Модель: $a(x) = w_1x + w_0$
- Два параметра w_1, w_0
- w_1 - тангенс угла наклона
- w_0 - где прямая пересекает ось ординат
- Функционал:

$$Q(w_0, w_1) = \frac{1}{l} \sum_{i=1}^l (w_1x_i + w_0 - y_i)^2$$

Парная регрессия



Выборка



Функционал ошибки

<https://spin.atomicobject.com/2014/06/24/gradient-descent-linear-regression/>

Парная регрессия

$$Q(w_0, w_1) = \frac{1}{l} \sum_{i=1}^l (w_1 x_i + w_0 - y_i)^2$$

$$\frac{\partial Q}{\partial w_1} = \frac{2}{l} \sum_{i=1}^l x_i (w_1 x_i + w_0 - y_i)$$

$$\frac{\partial Q}{\partial w_0} = \frac{2}{l} \sum_{i=1}^l (w_1 x_i + w_0 - y_i)$$

$$\nabla Q(w) = \left(\frac{2}{l} \sum_{i=1}^l x_i (w_1 x_i + w_0 - y_i), \frac{2}{l} \sum_{i=1}^l (w_1 x_i + w_0 - y_i) \right)$$

Линейная регрессия

$$Q(w) = \frac{1}{l} \sum_{i=1}^l (< w, x > - y_i)^2$$

$$\frac{\partial Q}{\partial w_1} = \frac{2}{l} \sum_{i=1}^l x_{i1} (< w, x > - y_i)$$

...

$$\frac{\partial Q}{\partial w_d} = \frac{2}{l} \sum_{i=1}^l x_{id} (< w, x > - y_i)$$

$$\nabla Q(w) = \frac{2}{l} X^T (Xw - y)$$

Начальное приближение

w_0 - инициализация весов

Например, из стандартного нормального распределения

Градиентный спуск

Повторять до сходимости:

$$w^t = w^{t-1} - \eta \nabla Q(w^{t-1})$$

Новая точка

Длина шага

Градиент в предыдущей точке

The diagram illustrates the components of the gradient descent update formula. Three light blue arrows point from descriptive labels to parts of the equation $w^t = w^{t-1} - \eta \nabla Q(w^{t-1})$. One arrow points from 'Новая точка' to w^t , another from 'Длина шага' to η , and a third from 'Градиент в предыдущей точке' to $\nabla Q(w^{t-1})$.

Длина шага

Повторять до сходимости:

$$w^t = w^{t-1} - \eta \nabla Q(w^{t-1})$$

Позволяет контролировать скорость обучения

Метод градиентного спуска

На каждом шаге (на каждой итерации метода) движемся в сторону антиградиента Функции потерь!

- Инициализируем веса $w_0^{(0)}, w_1^{(1)}, w_2^{(0)}, \dots, w_n^{(0)}$
- На каждом следующем шаге обновляем веса, сдвигаясь в направлении антиградиента функции потерь Q :

$$w_0^{(k)} = w_0^{(k-1)} - \nabla Q(w_0^{(k-1)}),$$

$$w_1^{(k)} = w_1^{(k-1)} - \nabla Q(w_1^{(k-1)}),$$

...

$$w_n^{(k)} = w_n^{(k-1)} - \nabla Q(w_n^{(k-1)}),$$

Сходимость

- Начальное приближение: w^0
- Повторять
$$w^t = w^{t-1} - \eta \nabla Q(w^{t-1})$$
- Останавливаем если
$$\|w^t - w^{t-1}\| < \varepsilon$$

Градиентный спуск

Останавливаем процесс, если

$$\|w^t - w^{t-1}\| < \varepsilon$$

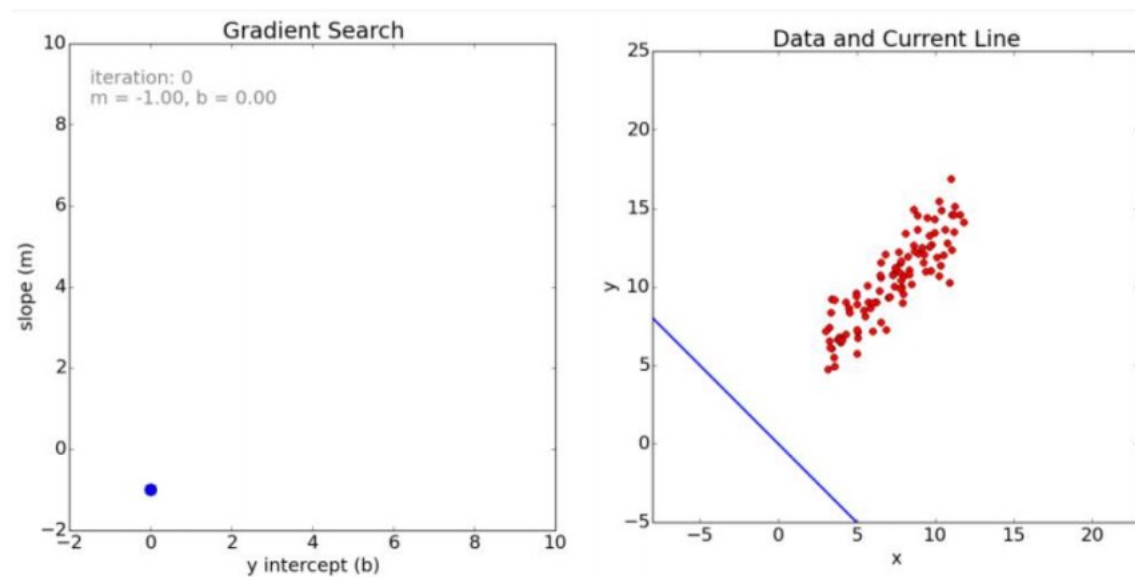
Другой вариант;

$$\|Q(w^t) - Q(w^{t-1})\| < \varepsilon$$

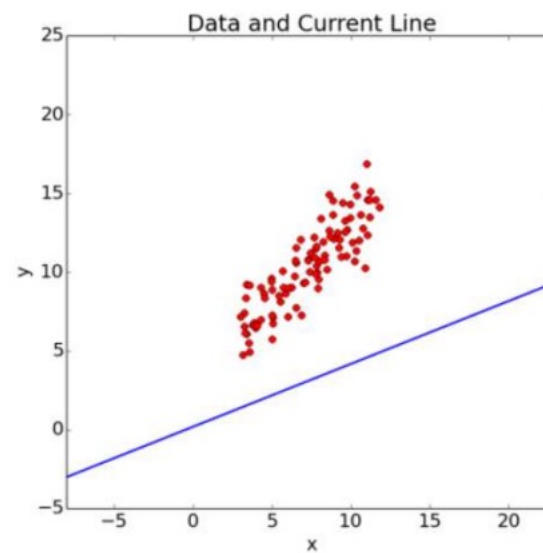
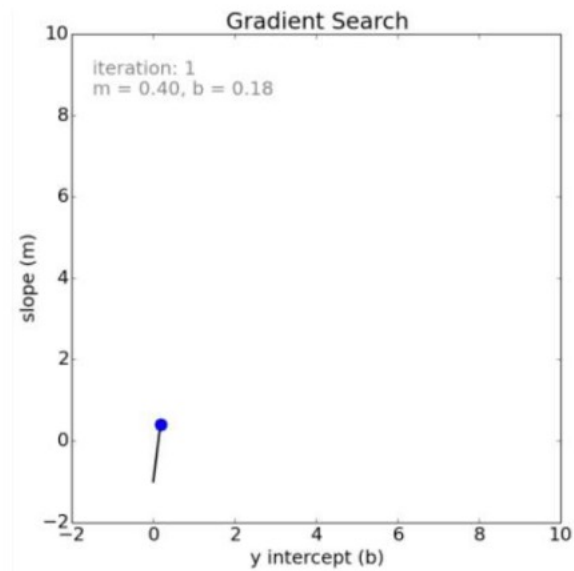
Другой вариант:

$$\|\nabla Q(w^t)\| < \varepsilon$$

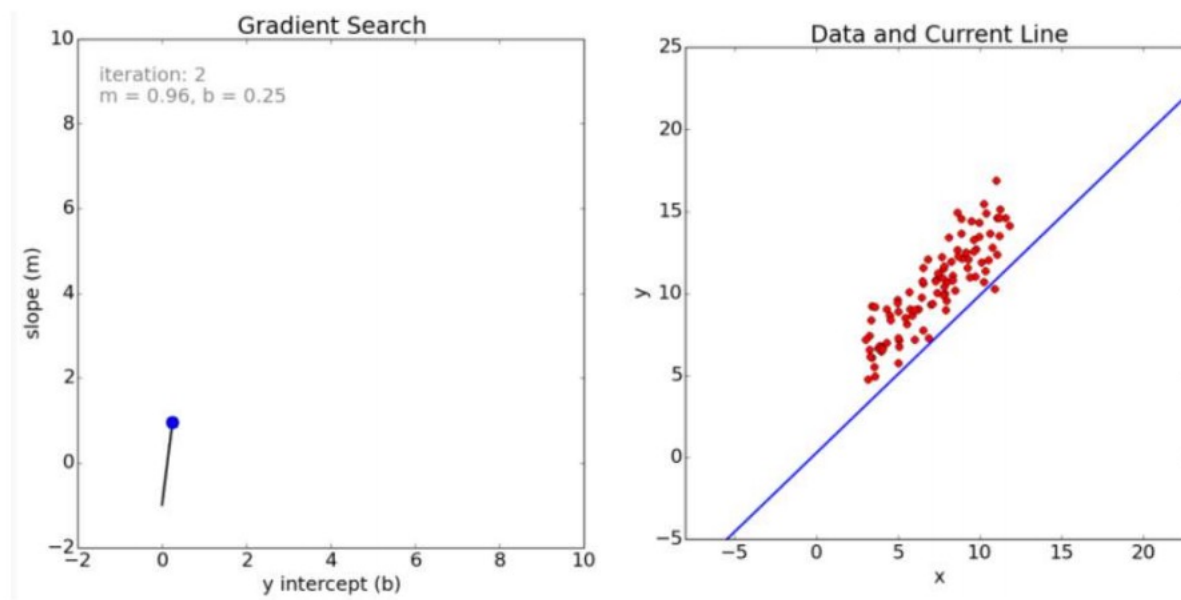
Парная регрессия



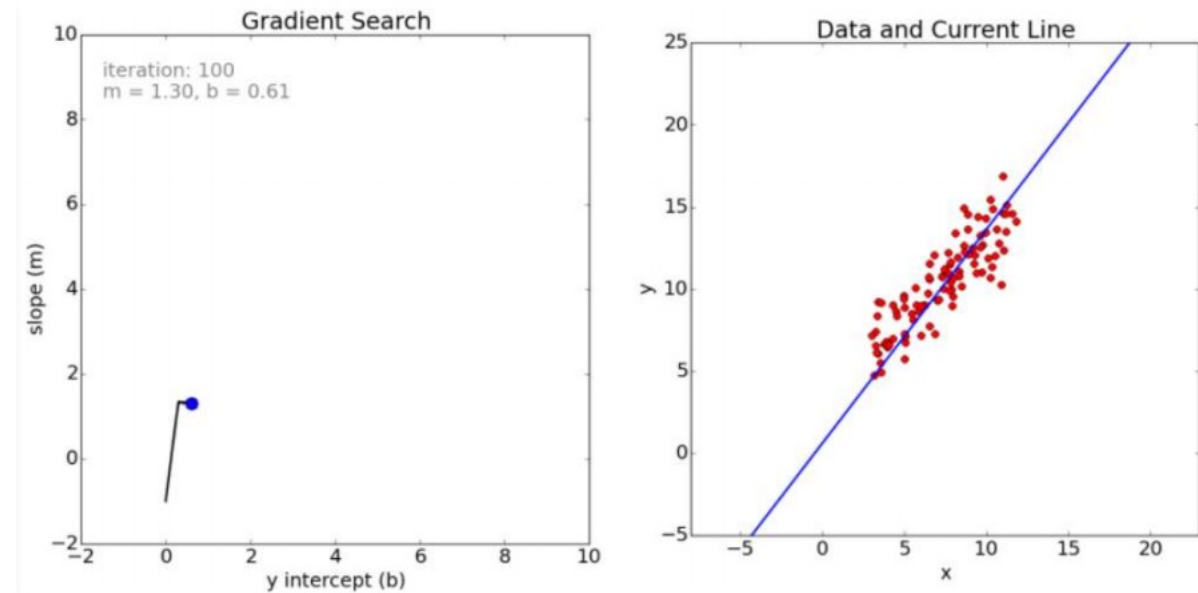
Парная регрессия



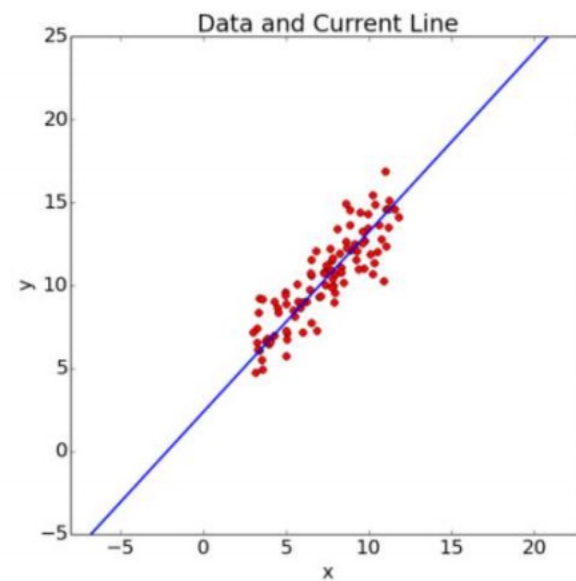
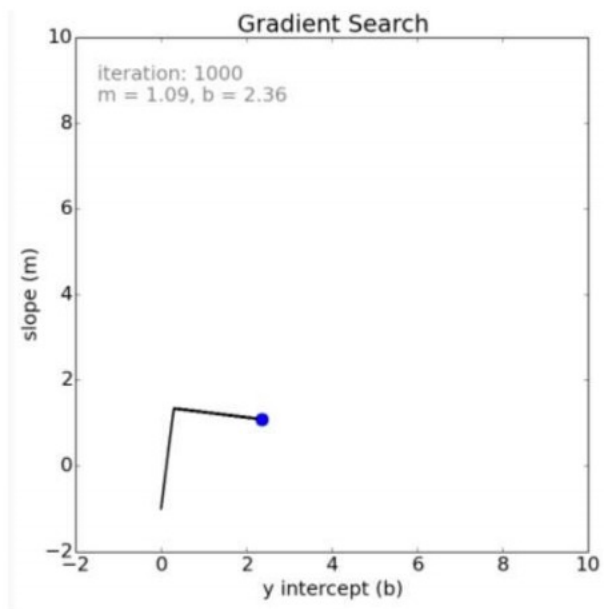
Парная регрессия



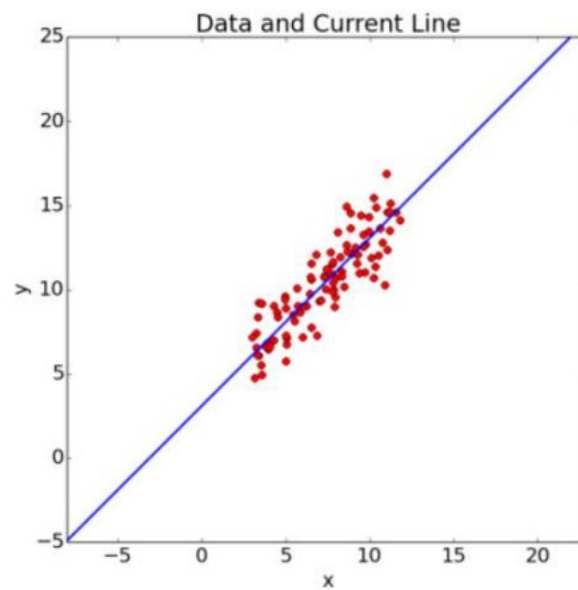
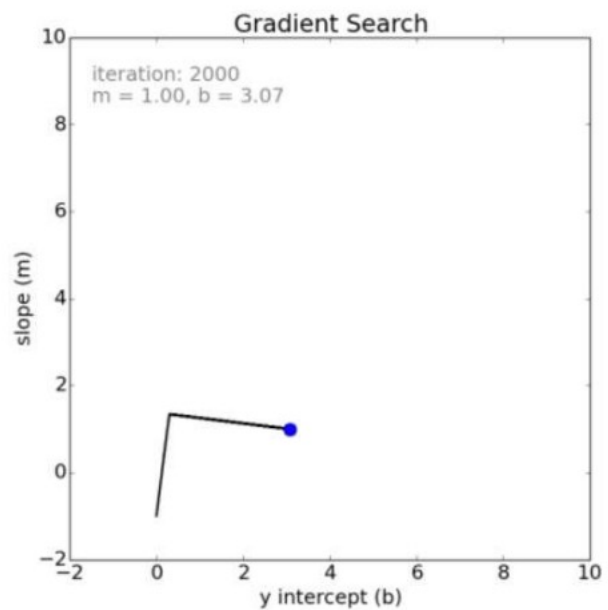
Парная регрессия



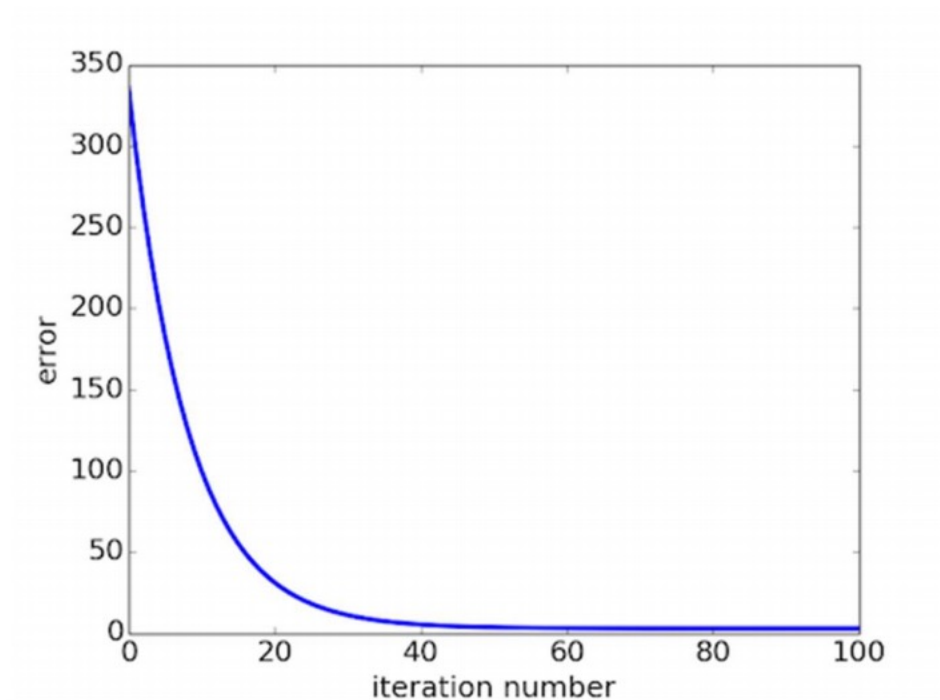
Парная регрессия



Парная регрессия

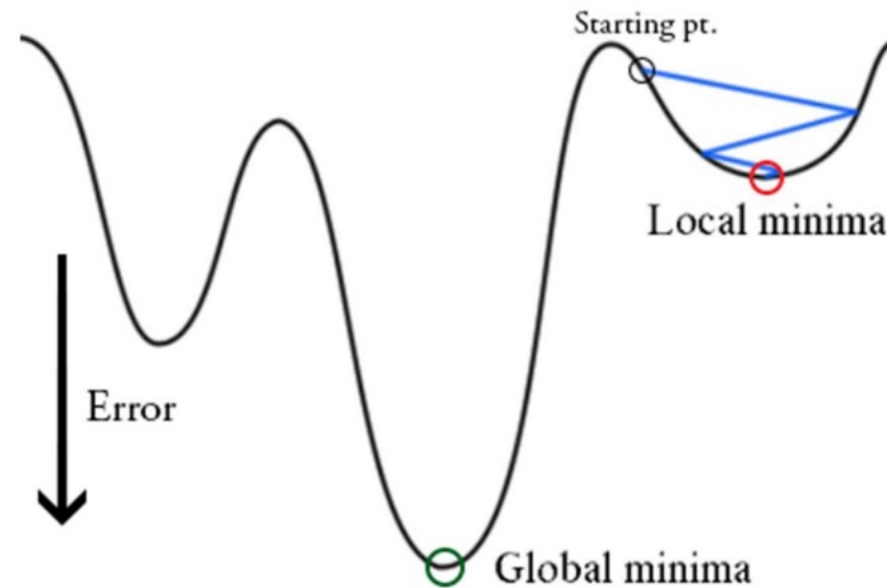


Функционал ошибки

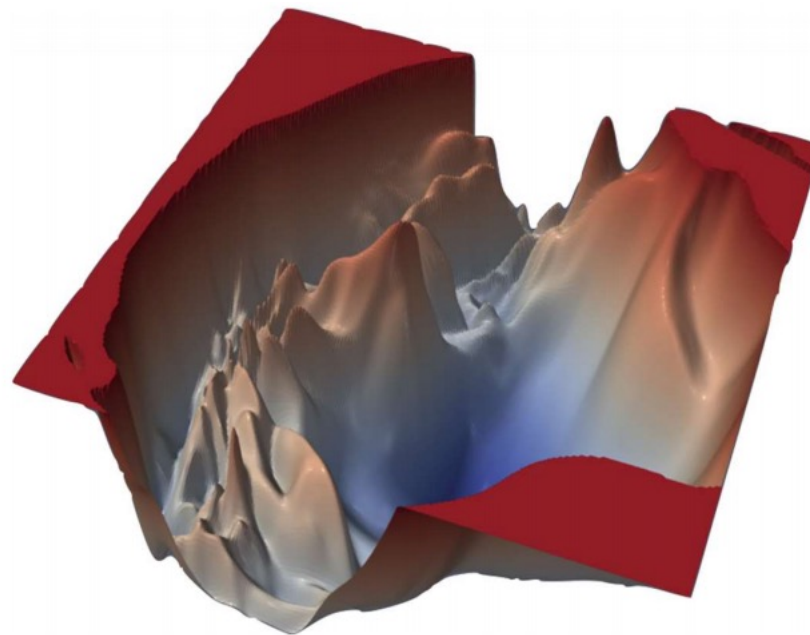
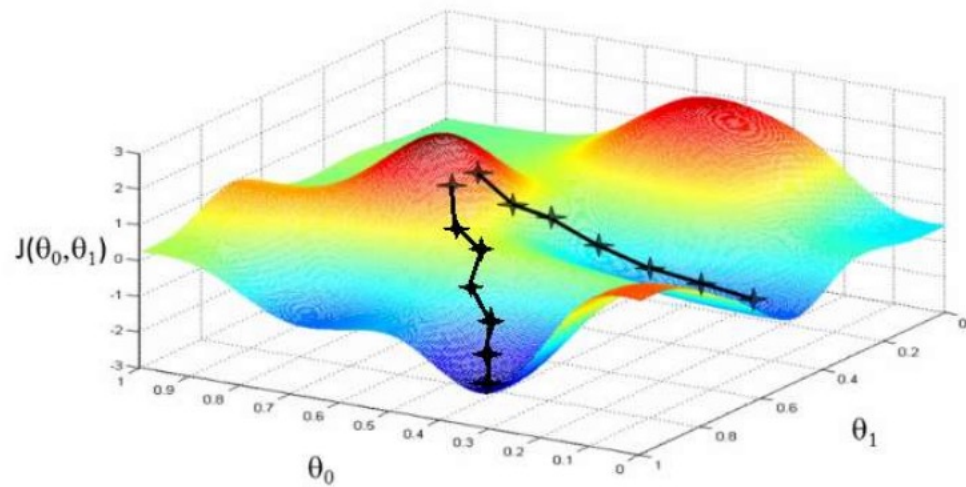


Локальные минимумы

Градиентный спуск находит только локальные минимумы



Локальные минимумы

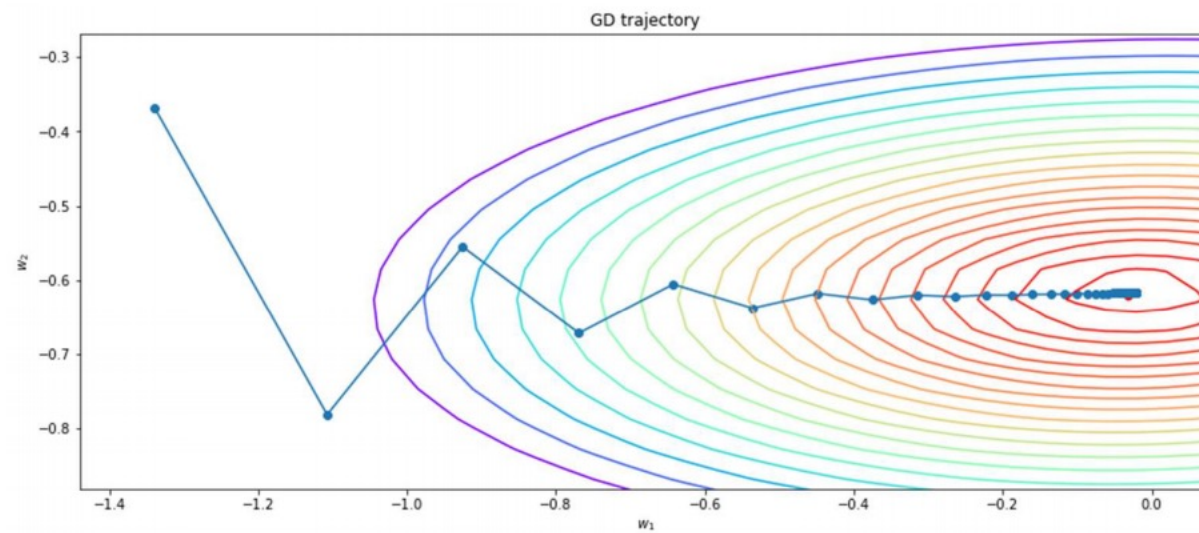


Длина шага

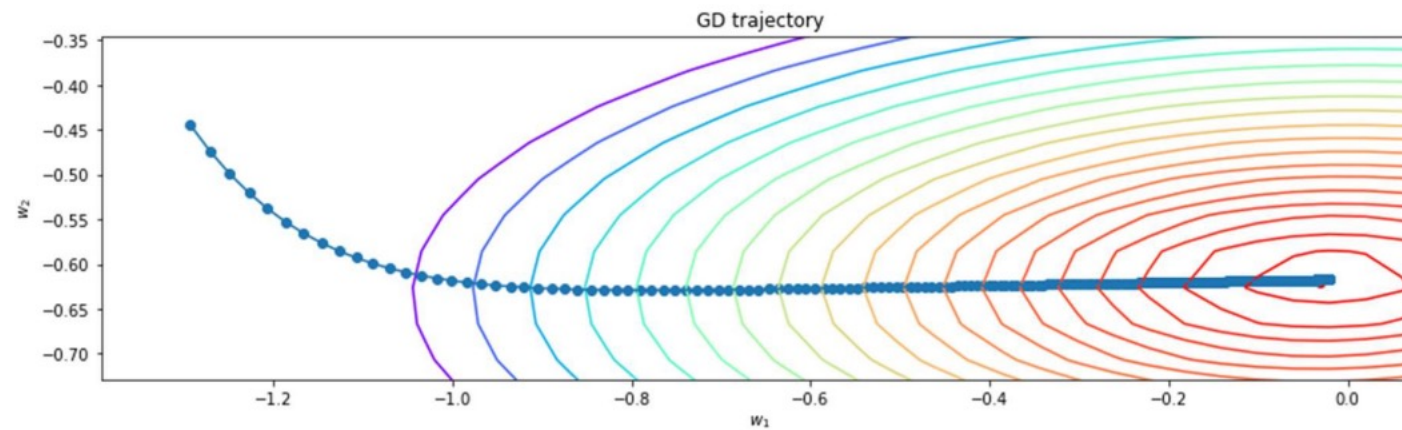
$$w^t = w^{t-1} - \eta \nabla Q(w^{t-1})$$

- Позволяет контролировать скорость обучения
- Если сделать длину шага недостаточно маленькой, градиентный спуск может разойтись
- Длина шага – гиперпараметр, который нужно подбирать

Длина шага



Длина шага



Переменная длина шага

$$w^t = w^{t-1} - \eta_t \nabla Q(w^{t-1})$$

- Длину шага можно менять в зависимости от шага
- Например: $\eta_t = \frac{1}{t}$
- Шаг наискорейшего спуска:
 $\eta_t = \operatorname{argmin}_\eta Q(w^t) = \operatorname{argmin}_\eta Q(w^{t-1} - \eta \nabla Q(w^{t-1}))$

Градиентный спуск

1. Начальное приближение : w^0

2. Повторять:

$$w^t = w^{t-1} - \eta_t \nabla Q(w^{t-1})$$

3. Останавливаемся, если

$$\|w^t - w^{t-1}\| < \varepsilon$$

Линейная регрессия

$$Q(w) = \frac{1}{l} \sum_{i=1}^l (< w, x > - y_i)^2$$

$$\frac{\partial Q}{\partial w_1} = \frac{2}{l} \sum_{i=1}^l x_{i1} (< w, x > - y_i)$$

...

$$\frac{\partial Q}{\partial w_d} = \frac{2}{l} \sum_{i=1}^l x_{id} (< w, x > - y_i)$$

$$\nabla Q(w) = \frac{2}{l} X^T (Xw - y)$$

Сложности градиентного спуска

- Для вычисления градиента, как правило, надо просуммировать что-то по всем объектам
- И это для одного маленького шага!

Оценка градиента

$$Q(w) = \frac{1}{l} \sum_{i=1}^l L(y_i, a(x_i))$$

- Градиент:

$$\nabla Q(w) = \frac{1}{l} \sum_{i=1}^l \nabla L(y_i, a(x_i))$$

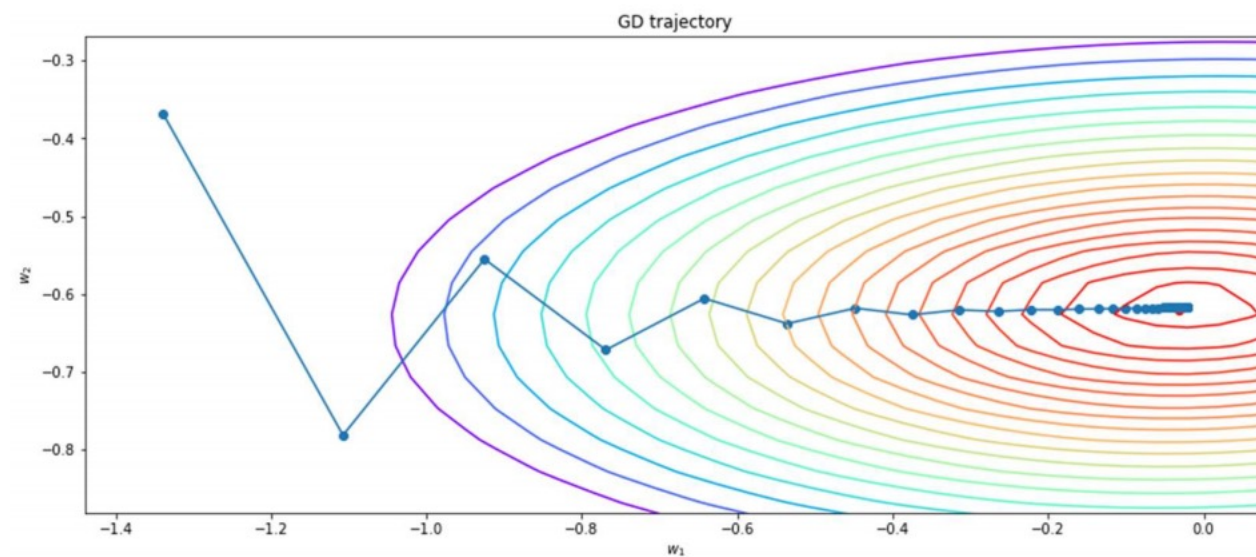
- Может, оценить градиент одним слагаемым?

$$\nabla Q(w) \approx \frac{1}{l} \sum_{i=1}^l \nabla L(y_i, a(x_i))$$

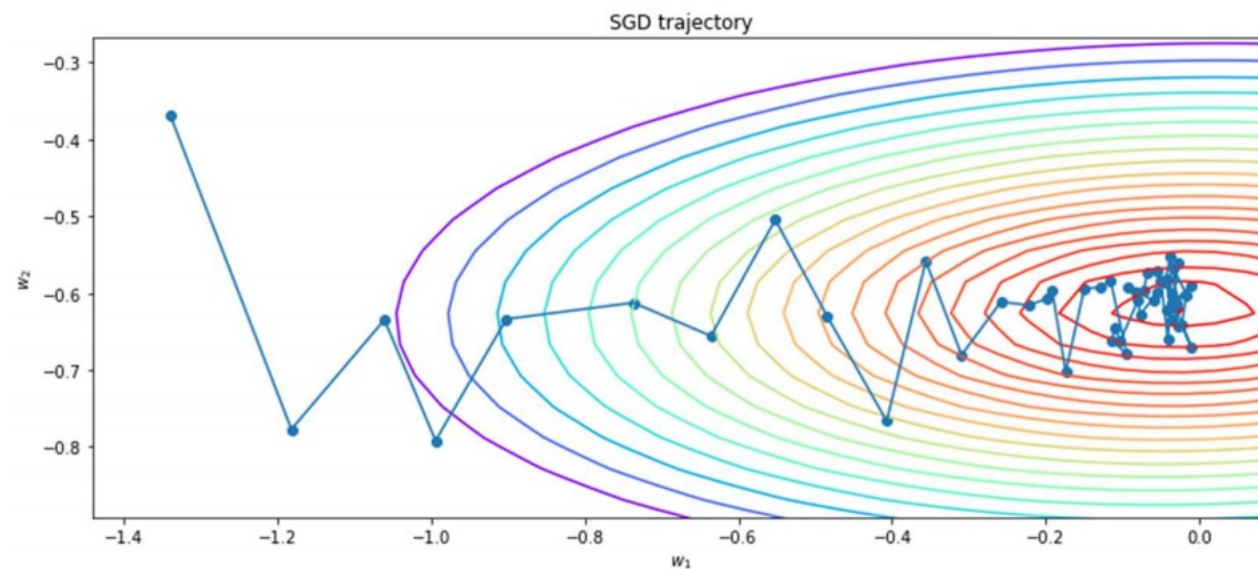
Стохастический градиентный спуск

1. Начальное приближение: w^0
2. Повторять, каждый раз выбирая случайный объект i_t : $w^t = w^{t-1} - \eta \nabla L(y_{i_t}, a(x_{i_t}))$
3. Останавливаемся, если $\|w^t - w^{t-1}\| < \varepsilon$

Градиентный спуск



Стохастический градиентный спуск

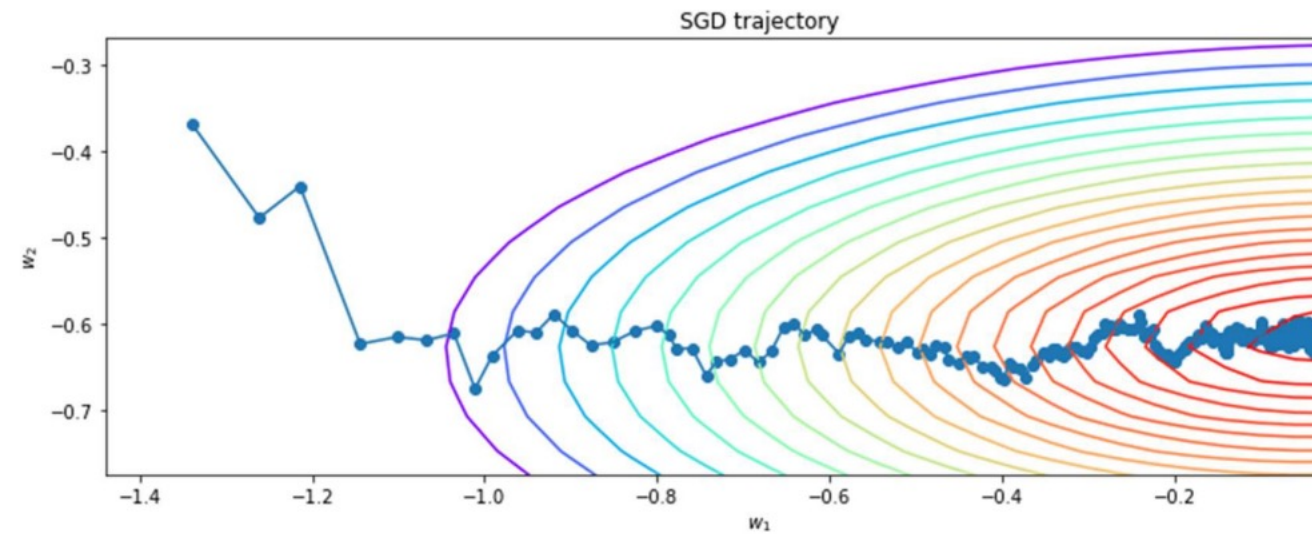


Стохастический градиентный спуск

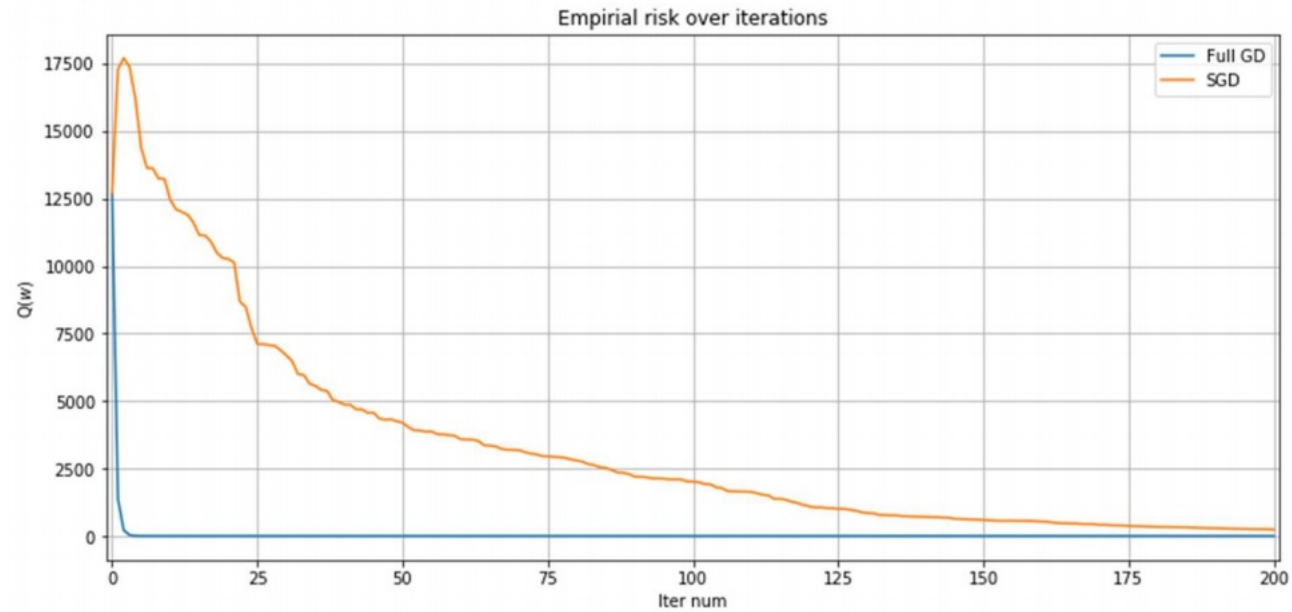
1. Начальное приближение: w^0
2. Повторять, каждый раз выбирая случайный объект i_t : $w^t = w^{t-1} - \eta_t \nabla L(y_{i_t}, a(x_{i_t}))$
3. Останавливаемся, если
$$\|w^t - w^{t-1}\| < \varepsilon$$

Стохастический градиентный спуск

$$\eta_t = \frac{0.1}{t^{0.3}}$$



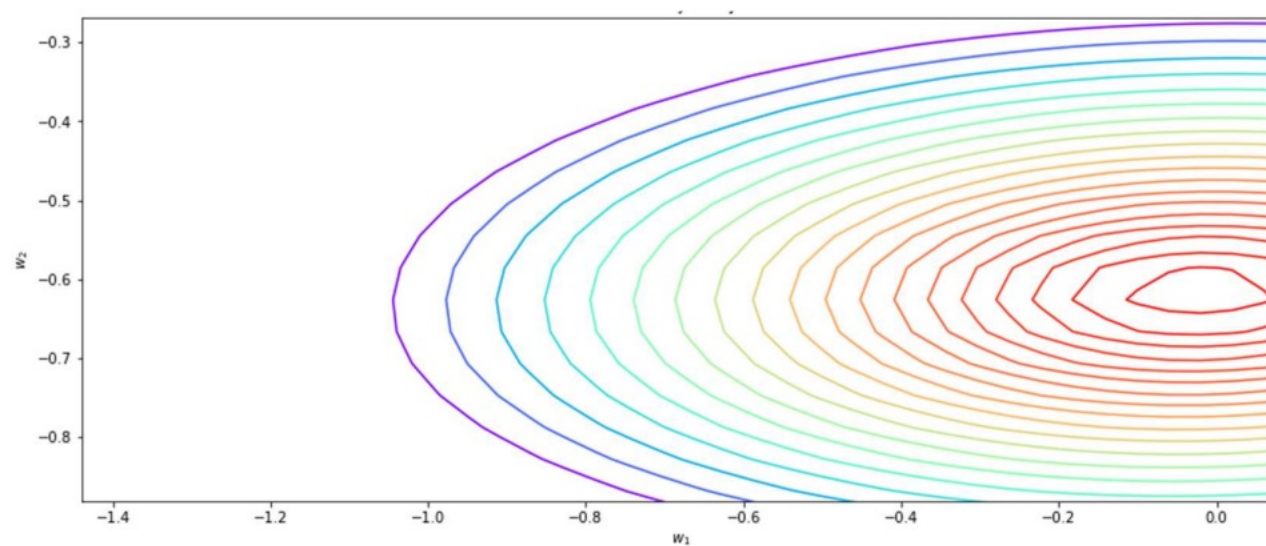
Стохастический градиентный спуск



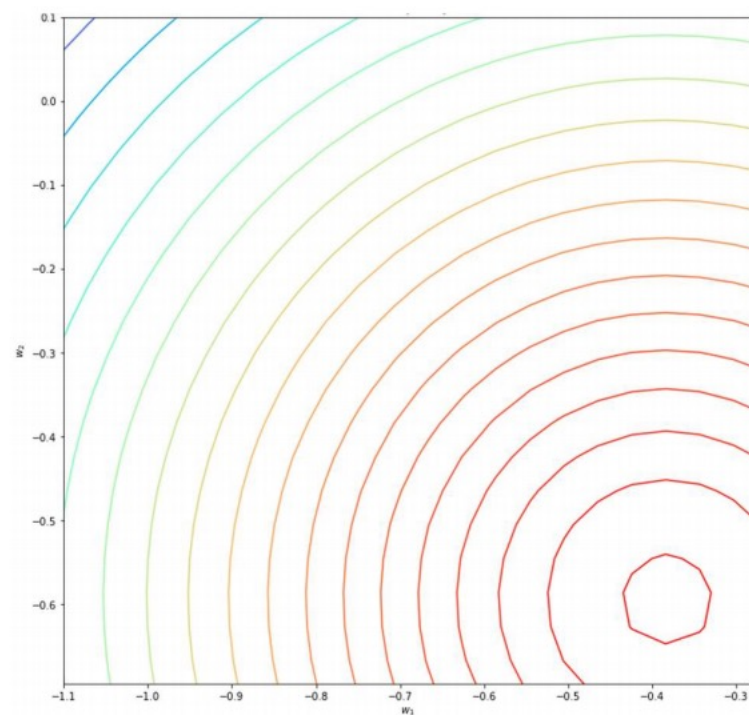
Стохастический градиентный спуск

1. Начальное приближение: w^0
2. Повторять, каждый раз выбирая m случайный объект i_1, \dots, i_m :
$$w^t = w^{t-1} - \eta_t \frac{1}{m} \sum_{j=1}^m \nabla L(y_{i_j}, a(x_{i_j}))$$
3. Останавливаемся, если
$$\|w^t - w^{t-1}\| < \varepsilon$$

Масштабирование признаков

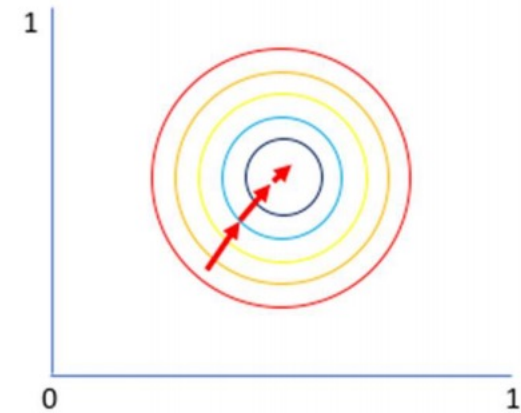
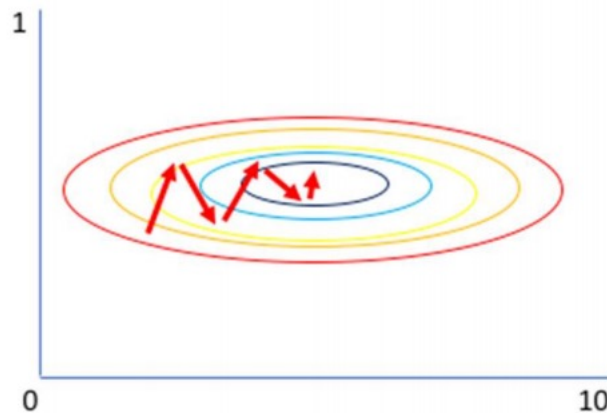


Масштабирование признаков



Масштабирование признаков

- До масштабирования значения градиента, соответствующие большим признакам, преобладают над остальными
- После масштабирования все параметры обновляются в равных пропорциях



Спасибо за внимание!



Ildar Safilo

@Ildar_Saf

irsafilo@gmail.com

<https://www.linkedin.com/in/isafilo/>